

## **17. User Interface Internationalization**

### **17.1 Operating With Non-US Languages**

Internationalization is the process of generalizing software so it can handle multiple languages (i.e., locales) and cultural conventions without the need for re-design or re-compilation. If an application designed for a US audience will be used in combined or coalition warfare operations, it needs to provide a user interface that matches users' expectations, interacts with users in their native language, and displays data in a manner that is consistent with users' cultural conventions. The guidelines in this section are provided to assist developers with a requirement to provide internationalized software and are not considered in determining COE style compliance.

#### **17.1.1 Character Rendering**

Languages can be categorized in terms of the characters or symbols in which they are written. To facilitate computer processing, a character set is defined for each language to contain its written letters, numbers, and punctuation marks, with each character in the set represented by a binary value. Most European languages, including English, are based on the Roman (or Latin) alphabet. Because these languages contain fewer than 200 basic characters (i.e., the 26 letters in the alphabet, with upper case, lower case, and accented variations), their character sets can be encoded in a single-byte. ASCII (American Standard Code for Information Interchange) is the most commonly used single-byte code set for representing English-language text for an American user base. However, because ASCII encodes a limited number of characters, it is insufficient for representing text in languages other than American English.

While most European languages are based on the Roman alphabet, many of them contain extended characters (i.e., ones that do not exist in English and are not available in ASCII) in their character sets. These characters include accented vowels such as é and ê; characters such as the French ç, the Spanish ñ, and the German ß and ü; and combined characters such as æ. In addition, some European languages may not use the entire Roman alphabet; Italian, for example, lacks the letter k. Despite these variations, all text in Roman-based languages is written from left to right, with each new character appended to the right of the previous character. Furthermore, the appearance of a character and its order within a character sequence do not change as new characters are entered.

While languages with fewer than 200 characters can be encoded in a single byte, complex languages such as Chinese, Japanese, and Korean contain several thousand unique ideographic symbols from which words are composed. Encoding such a large character set requires two (or more) bytes per character rather than one. Unicode provides the capability to represent multi-byte character sets and is the preferred approach for encoding the alphabets, ideograph sets, and symbols for all languages in the world. In multi-byte languages, there may be no "natural order" to characters vis-à-vis sorting, and no concept of or distinction between upper-case and lower-case forms. Because

characters are composed of many strokes, text in these languages may require more space to display, and text entry may be a more complex process than in single-byte languages

In most Roman-based languages, each character is rendered as a separate symbol of fixed shape, and characters are written in the approximate order in which they are pronounced. In other languages, however, the way in which characters are rendered graphically depends on their linguistic context. In Arabic, for example, a character may be displayed in several different forms, depending on whether it is displayed alone or as the first, middle, or last character in a word. In contextual languages, the characters that make up a symbol may be entered in several different orders, and entering a new symbol may change or even eliminate a previously entered symbol. As a result, complex algorithms may be needed to manage text line length (e.g., line breaks, justification), support editing individual characters within a symbol, and provide search and sort features that can recognize multiple encodings of the same symbol.

Most languages are unidirectional; i.e., lines of text are presented uniformly from left to right or from top to bottom. Although Asian languages such as Chinese, Japanese, and Korean may present horizontal and vertical text on the same page, they are considered to be unidirectional because they do not mix directions in a single line of text. By contrast, Arabic and Hebrew are bi-directional; text in these languages is written from right to left, but numbers and foreign words in the same text are written from left to right. Because the direction of text entry may change from one character to the next, appropriate text handling procedures must be available in both right-to-left and left-to-right text.

### **17.1.2 Structural Rules for Character Handling**

An application is able to accept and process all of the characters in the character set used by the target language (i.e., the one to which the application is being converted). Because languages differ in their structural rules for character handling, assumptions made when processing a US character set may be inappropriate or inaccurate when applied to a language with extended characters. In particular, a US application is likely to require modification in order to correctly handle case conversion, ligatures, special characters, and word and character boundaries in the target language.

Case conversion. In US software, case conversion is usually performed by adding or subtracting a constant (i.e., 32) to or from the ASCII code for the character. In extended character sets, case conversion is more complicated because there is no constant difference between the numerical equivalents for upper- and lower-case representations of characters. In addition, the distinction made in English between upper-case and lower-case letters may be ambiguous or not exist at all in other languages. For example, Chinese, Japanese, and Korean have no case distinction. In other languages, an accented vowel in lower case may retain its accent in upper case, or the accent may disappear; alternatively, a vowel in upper case may or may not contain an accent in lower case, depending on the word and the language.

Ligatures. Ligatures are sequences of characters that are treated as a unit; e.g., æ is a combination of a and e. In some languages, ligatures can be entered as a single character or as two separate characters. In the latter case, both letters would be capitalized in words that are proper nouns; for example, Iceland is written as IJsland in Dutch. Ligatures occur frequently in contextual languages such as Arabic. A US application may require revision in order to handle any ligatures that occur in the target language.

Special characters. Because languages differ in the meaning assigned to special characters, a US application that uses characters such as apostrophes as delimiters in a text string and restricts their use to this function may require modification when converted to certain European languages. For example, French and Italian replace the terminal vowel in an article by an apostrophe when the following noun has an initial vowel. In addition, some languages include special characters that may not be present in English or use these characters in ways that differ from US usage. For example, Spanish starts exclamations and questions with inverted exclamation mark and question mark characters, while French includes a space between the last word of a sentence and a concluding exclamation mark or question mark. Finally, diacritical marks (i.e., the signs modifying the value or sound of characters) may have different meanings in different languages. For example, certain diacritical marks specify the doubling of consonants in Arabic but may indicate pitch in Vietnamese.

Word and character boundaries. A word consists of a string of characters preceded and followed by delimiters. In a US application, these delimiters are usually blanks or spaces but may also include the unused portion of the Roman character set. This latter approach can be problematic when converting to a Roman-based language where these characters need to be interpreted as part of the word and not considered as delimiters. In addition, in some languages, a blank is acceptable as a numeric or phrase separator and so would not be appropriate to use as a standard word delimiter.

## **17.2 Text Translation**

### **17.2.1 Creating Internationalized English Text**

The process of translating the text displayed by a US application begins with the creation of an “internationalized English” version of the text. All of the text is reviewed and, if necessary, modified to ensure that it is easy to understand and use. Message text (e.g., in message windows, online help) is presented in short, simple, declarative sentences whenever possible. Excessive use of subordinate and coordinating phrases is avoided, and ideas are expressed as concretely as possible. Ambiguous language, humor, jargon, and cryptic messages are likely to cause difficulty for non-US audiences and so need to be eliminated. Likewise, compound adjectives, strings of nouns, long sentences with many ideas, and negative questions can be difficult to understand and so are not used. If an application needs to explain a series of concepts, they are presented in the form of a list, rather than in a text string separated by commas.

The content of each window in an application is checked for US-specific language prior to translation to minimize the likelihood of misinterpretation by the target audience. The goal is to use only those terms that are employed in the same way throughout the English-speaking world. The use of acronyms and abbreviations is limited since many are not recognized internationally and may have different meaning, depending on where they are used. Similarly, when large numbers are presented, they are written as numerals; for example, the term “billion” means one thousand million in the US but one million million in some European countries. To minimize confusion, the names of months are written out when they appear as part of a date. For example, 06/10/94 can be interpreted as June 10 or October 6, depending on the user’s cultural background.

An application avoids presenting examples that may be uniquely American. A generic term is used, rather than what something is called in the US. For example, “stock exchange index” is an international term while “the Dow” is specifically American and “the FTSE 100” is specifically British. The messages in an application are reviewed to determine if they may be interpreted by the target audience in ways other than intended. For example, “as soon as possible” means “immediately” in the US but “when convenient” in other countries. Terms such as “left hand” can be offensive in some cultures and need to be replaced with “on the left” or “left side.”

### **17.2.2 Translating Text and Messages**

Once an internationalized English version of the text displayed by an application is created, it is translated into the target language. If appropriate, the translation is tailored to the target language in the specific country or region that will be using the application. Translated text contains proper technical terminology, especially when the terminology may differ from conversational expression in the target language. Care is taken to maintain distinctions in terminology that may be translated into the same text string in the target language. For example, Cancel and Undo are normally translated as “annulez” in French, even though these commands have different meanings in English. Similarly,

spelling and grammar can differ among varieties of a single language and so need to be adapted accordingly.

The same terminology is used in both the user interface and documentation for an application. While the goal is to provide an accurate translation of all text into the target language, it is acceptable to use US words in the text if the target language does not have an adequate vocabulary of technical words or if the target audience is accustomed to dealing with US terminology.

The accuracy of translated messages is verified since it is possible for the text of one message to be the same as another, especially if the original messages were very similar or were worded ambiguously. In addition, the meaningfulness of the translation is checked against the situation that invoked it to ensure that the information being conveyed in the original message is also conveyed in the translated version.

Translated text is reviewed to ensure that it makes grammatical sense. Some applications construct messages from two or more substrings; for example, the name of the file being deleted is inserted into the text string requesting the user to confirm deletion of the file. While this approach may work in English, the linguistic characteristics (e.g., gender, word order, special characters) of the resulting text can be awkward or inappropriate when translated into the target language. When messages are constructed by nesting or concatenating strings, the translated text that results is frequently meaningless or syntactically impossible because words or phrases were not modified to fit the grammatical rules of the target language.

Translated text uses the same character set and font as the rest of the application. Line breaks and other format changes that may have been introduced with the translation are checked for accuracy since US rules for hyphenation, punctuation, or capitalization are likely to be different from those in the target language. If typographic variations such as italics or boldface have been added as part of the translation, they are checked to ensure that they are suitable in the target language. Language environments have evolved unique rules defining how elements such as title lines, bulleted lists, and footnotes are used to distinguish among levels of expression and to indicate how expressions are related. The appearance of the translated text is adapted as needed to satisfy these rules.

### **17.2.3 Translating Documentation**

The documentation for an internationalized application describes a representative sample of the internationalized capabilities provided by the software. For example, an explanation of how a sorting function works describes the kinds of sorts that are performed, explains that the current locale affects the output, and provides several examples that are representative of the locales supported.

Documentation text contains simplified English. Whenever possible, a single term is selected to express a concept, and the use of synonyms is minimized. However, these changes are made in such a way as to not reduce the precision of the text, create awkward

phrasing (with an increase in overall text length), or produce unacceptably dull or boring text. References (e.g., to sample users) and examples that are specific to one culture are modified to be more international in focus. Any graphic symbols used in the documentation are reviewed to minimize the extent to which they are culture-specific. If necessary, a table is provided that lists the symbols and their interpretation. Finally, documentation sections such as glossaries and indexes are expanded as needed to help non-US readers find information. For example, glossaries define words that may have a different technical meaning or not exist at all in the locales supported by an application.

Because the order of items in a sorted list usually changes following translation, references to the position of items in the list are removed from documentation. Similarly, when collation sequences are described, the results are not described as sorted “alphabetically” since ideographic languages cannot be sorted alphabetically. Instead, output is described as appearing “in sorted order as determined by the current locale.” Other changes needed when internationalizing documentation include ensuring that terms such as ASCII, text, byte, and character are used appropriately, describing any assumptions made about date and time formats, replacing references to Yes and No (e.g., when describing actions in response to a message) with words that are appropriate to the locale, and presenting the names of any individuals (e.g., sample users) in an order that is correct for the specific culture. Finally, lengthy text explanations (e.g., in online help, training materials, or other documentation) may need to be restructured or reorganized so they follow the rules used by the target audience in organizing technical discussions or sequences of explanations.

## 17.3 Text Input Methods

### 17.3.1 Keyboards and Keyboard Input

A workstation usually has a keyboard layout tailored to the target language, and this layout may be different from the one available on US workstations. Conventions concerning the location of characters vary from language to language and sometimes from country to country within the same language. For example, the German keyboard reverses Z and Y from their positions on the US keyboard, and the Spanish keyboard has a different layout in Spain than in Latin America. In addition, languages may add, omit, or change the characters on a keyboard. For example, the British keyboard contains the currency symbol for pound instead of #, and the Spanish keyboard has ñ where the US keyboard has L. Non-US keyboards may mark each key with up to four different characters. Users press modifier keys (e.g., SHIFT and/or ALT) in combination with the key to enter the various characters marked on the key.

Because computers respond to specific physical keypresses regardless of what markings appear on the keys, a different keyboard may not be required when converting a US application into another language. A keyboard can be adapted by replacing the symbols on each key, either with adhesive labels or new key covers. An application then maps the individual keystrokes to the character set for the other language and displays the appropriate characters. If this approach is used, the function keys on the keyboard need to be mapped to the same actions as in the original software, and any messages generated when these keys are pressed are displayed as they were prior to the conversion.

Languages where diacritical marks are used extensively (e.g., French) usually provide keyboards that allow users to generate characters with these marks with a single keystroke. However, because English has very few accents, users with a US keyboard have to execute a combination of keystrokes in order to enter an extended character. An application can use “dead” keys or a compose-based method to produce this type of input.

With “dead” keys, the keystrokes consist of a “dead” (i.e., nonspacing) key, followed by the character (e.g., a vowel) to be displayed with an accent. A different dead key is assigned to each accent. When a dead key is pressed, a text input mode is invoked; the symbol on the key is not displayed, and the text cursor does not move. The mode is automatically disabled following the next keystroke; the appropriate dead key is pressed each time an accented character is entered. If an invalid character (e.g., a consonant) is entered, the character is displayed without an accent, and feedback (e.g., a beep) is provided to indicate that the keystroke was invalid.

In a compose-based input method, when a predefined control key is pressed, a text input mode is invoked that forms the next two keystrokes into a single character. When the first character (e.g., a vowel) is typed, nothing is displayed by an application. When the second character (e.g., the diacritic) is entered, the completed character is displayed, and the input mode is automatically exited.

### **17.3.2 Approaches to Text Entry**

Text entry using pre-edit methods. In most languages, users perform text entry by typing directly into a text box. However, if a keyboard cannot produce all of the symbols in a target language, a pre-edit step may be needed. Users type characters from the keyboard, usually into a pre-edit area, and then execute an action to convert the characters into other symbols appropriate to the language. These symbols are then displayed in the text box.

When a pre-edit step is required, text entry can be performed on-the-spot, over-the-spot, or off-the-spot. On-the-spot means that as users type, the characters appear directly in the text box which can contain both text in unconverted form and converted symbols. Although more difficult to implement, this approach is preferred because it is more similar to text entry as normally performed by users. In over-the-spot, a separate pre-edit area is provided for each text box; when users convert their input into final form, the symbols are displayed in the appropriate text box. Off-the-spot also provides a separate pre-edit area but uses the same area for multiple text boxes; in this case, when users convert their input into final form, the symbols are displayed in the text box that has input focus.

When text entry includes a pre-edit step, an application provides feedback concerning the status of the input after users enter text in the pre-edit area and then execute an action to convert the input into final form. If insufficient information is available to perform the conversion, an application can prompt users to enter more pre-edit text, present them with a list of choices from which to select, or indicate that the conversion has failed. If an on-the-spot approach is implemented, the text box provides a visual distinction (e.g., a different text font or color) between original input and converted text so users can easily distinguish between the two. If the pre-edit area is provided in a separate dialog window, the window is modeless so users are not restricted to only performing text entry.

Text entry in languages with large character sets. Several options are available to support keyboard input in languages that have large character sets. Whatever method is selected must be able to accommodate context-specific variations within the language as users perform text entry. With each keystroke, converted text changes as needed in order to create a new compound character or add a mark to a previous character.

With the first option, the component elements of each character are marked on the keyboard. As users press individual keys, the elements are displayed. When a character is complete, it is displayed in place of its components. This method has been used to perform text entry in Chinese and Korean.

With the second option, users enter each character phonetically, and the phonetic form is automatically translated into the correct character. When more than one character has the same pronunciation, users are presented with a list of phonetically similar characters from which to choose. For example, users enter a root or radical character from the keyboard, then select additional strokes to complete the character from a set displayed by the

application. This method has been used to convert Roman characters to Chinese ideographs, and Hiragana and Katakana characters to Japanese Kanji.

With the third option, users type the decimal or hexadecimal encoded value for a character or select the value from a list. If the value matches an entry in the code set, the corresponding character is displayed by the application.

Text entry in mixed character sets. Users may need to perform text entry in more than one character set (e.g., English and Korean) or in multiple locale-specific character sets (e.g., Kanji and Katakana). This flexibility can be provided by defining text input modes in each character set, along with a special keyboard character that allows users to toggle between the character sets as desired. Users with a keyboard where two character sets are marked on the keys select one of the modes to begin text entry. All of the typed text is interpreted in this character set. When the special character is encountered, the text mode toggles to the other character set and all subsequent input is interpreted in this set.

Text entry in bi-directional languages. Because bi-directional languages write text in both right-to-left and left-to-right directions, text entry may be performed in either direction, depending on the contents of a text box, and may require input in both directions within a single box. An application can provide automatic handling of directionality based on the characters being typed, or it can support an input mode that users invoke to switch the language and direction of text entry.

### **17.3.3 Other Text Entry Actions**

The text cursor remains visible during text entry to indicate the locus of typed input. In addition, the text cursor does not disappear from view as it moves from one character to the next in a string of single-byte and multi-byte characters. If the target language is bi-directional, an application can support multiple text cursors within a single text area, one indicating when text can be added in the current input direction and the other marking the last place where the direction of input changed. In contextual languages (i.e., where the appearance of existing text can change as new characters are entered), the text insertion point can move backward or forward as users perform text entry; in this case, the text cursor is displayed in a manner that is consistent with the movement of the text insertion point.

The arrow keys move the text cursor in the direction of the arrow regardless of the direction in which text is currently being entered. DELETE deletes text in the direction opposite to the direction in which text is being entered.

If an application is being converted to a contextual language, it needs to define how certain keystrokes affect a compound symbol that is composed of several separate characters. For example, an application needs to determine when DELETE cancels the previous keystroke (i.e., removes a character) or deletes the entire symbol. In addition, an application may need to limit the ability to insert or delete individual characters in a

word if these actions would change neighboring characters or alter the appearance of the word in unintended or confusing ways.

## 17.4 Internationalizing User Interface Features

### 17.4.1 Text Expansion

When text is translated into another language, the result is often longer than the original English. For example, the phrase “message pop-up” translates to “Nachrichtenüberlagerrungsfenster” in German and “janela de sobreposiçao de mensagem” in Portuguese. The increase in text length may be as much as 200 percent, depending on the length of the original text. Some of this increase may result from the addition of spaces that were not present in the original text. Table 17-1 lists allowances for expansion recommended by MS Windows based on text length in English. This table refers to the number of characters in a message, with characters in multi-byte languages (e.g., Japanese) taking two bytes per character.

**Table 17-1. Allowances for text expansion.**

| <u>Length of English Text</u> | <u>Additional Space Required</u> |
|-------------------------------|----------------------------------|
| Up to 10 characters           | 200 percent                      |
| 11 - 20 characters            | 100 percent                      |
| 21 - 30 characters            | 80 percent                       |
| 31 - 50 characters            | 60 percent                       |
| 51 - 70 characters            | 40 percent                       |
| Over 70 characters            | 30 percent                       |

Note: This table was taken from the Microsoft Windows Software Development Kit -- Additional Windows Development Notes, as published in Software Internationalization and Localization: An Introduction.

Translated text may require adjustments in the horizontal spacing between specific pairs of characters. For example, in English, the characters f and i look better when displayed closer together than other pairs of characters. The horizontal spacing algorithms used by an application need to accommodate adjustments in the spacing of non-US characters, including pairs of characters (e.g., æ) that may be part of an extended character set.

The height of a line of translated text may be twice the height of the text in English. Roman-based languages may supplement the character set with diacritical marks that extend above or below the basic symbol. In non-Roman languages, marks may be stacked two or three high, and small versions of characters may be placed above, below, or beside the primary symbols, causing wide variations in the height of each text line. Because of their complexity, ideographs require more space to display the strokes within them. For example, some complex Chinese characters may need to be displayed at least 50 percent larger than alphabetic characters in order to be readable. The minimum size for ideographs is usually 16 x 16 pixels. Translated text may also require adjustments to the vertical spacing between lines to ensure legibility and readability when displayed by

an application or printed. Extended characters, and in particular those with diacritical marks, may require additional spacing, especially when printed in upper case.

It is likely that the size and placement of controls in application windows will require adjustment in order to accommodate text expansion following translation. Menus and dialog windows will also need to increase in size in order to accommodate the longer text. Similarly, more vertical space may be needed in window components such as the title bar to accommodate larger character size, especially in languages such as Chinese, Japanese, and Korean. Finally, the size of the text included in the label of a window icon may need to increase to accommodate an extended character set, and the icon graphic may also contain embedded text that needs to be translated.

Each application window needs to be checked to ensure that all of the translated text fits properly within the window and that individual controls are positioned correctly within each window area. For example, when column headings are translated, they may be longer than the data they contain and need to be broken into more than one line of text. Similarly, when text labels are translated, the placement of the associated text boxes may be altered and require repositioning in order to be properly aligned within the window.

### **17.4.2 Nonlinguistic Text Features**

Capitalization, punctuation, and word order. Text displayed by an application follows the rules for capitalization, punctuation, and word order used in the target language. For example, in German, all nouns are capitalized, regardless of their position within a phrase or sentence. Depending on the language, quotation marks may be displayed as “quotation”, «quotation», »quotation«, or „quotation.“ Interrogatory sentences in Spanish begin with an inverted question mark and end with a question mark in normal orientation. Adjectives precede nouns in English word order but may follow nouns in other languages.

Hyphenation. An application performs hyphenation in a manner that is consistent with the rules of the target language. These rules may call for changing the characters in a word when it is hyphenated at the end of a line, or placing hyphens between individual words when they extend beyond the end of a line. For example, in German, “drucken” and “heißen” become “druk-ken” and “heis-sen” when hyphenated; in French, a hyphen is added between a personal pronoun and “même” (e.g., “eux-même”) when these words extend beyond the end of a line.

Justification. The justification routines used by an application conform to the rules of the target language and may require some character-processing logic in order to do so. For example, in Asian languages where spaces are not used to delimit words, line breaks can occur anywhere within a word. However, because symbols are represented by a multi-byte character, line breaks cannot occur within a symbol nor can punctuation be the first character on a new line. Alternatively, languages such as Arabic and Hindi do not allow breaks within words. In this case, the justification algorithm used by an application must

accommodate this restriction and be able to produce justified text without excessive space between words.

Abbreviations. The abbreviations used in US software may have different meanings in other languages or not be used at all. For example, while # is commonly used as an abbreviation for number, this character is not meaningful outside the US. The symbol @ means “at” in the US but “each” in the United Kingdom. The abbreviations for ordinals are 1st, 2nd, 3rd, etc. in the US, but 1<sup>o</sup>, 2<sup>o</sup>, 3<sup>o</sup> or 1<sup>a</sup>, 2<sup>a</sup>, 3<sup>a</sup> in other languages, depending on the gender of the subject.

Typography. If an application displays text in a Roman-based character set, it supports the fonts (e.g., Times and Helvetica), sizes (e.g., 10-point, 12-point), and styles (e.g., regular, italic, bold) that are normally available in the target language. An application also accommodates any unique typographic conventions when displaying translated text. For example, stress in writing is indicated through the use of italics in English but by letter spacing or boldface in European languages. In Japanese, stress is indicated by underlining characters, putting a light gray background behind them, or writing the text in Katakana.

Reordering sorted information. If an application presents sets of related items (e.g., in lists, option menus) in alphabetical order, it reorders the items after translation based on a sort sequence that is meaningful to the target audience. The most appropriate order may vary by application and depend on the information displayed in the items. Section 17.4.6 provides additional information on sorting and collation.

### **17.4.3 Data Formats**

An application is able to recognize and correctly handle the range of formats that are used to express data in the target language. The labels for all data entry areas are modified to include the appropriate unit of measurement. An application either converts the data format to one that is familiar to users or provides the capability to display data in alternate formats so users can select the one that is most meaningful to them. For example, US users prefer to measure length in feet and yards while European users are more familiar with the metric system. If the content of data entry areas is not converted to a format that is familiar to the target audience, then the data format is included as part of the data label. If an application supports converting to and from both millimeters and inches, the number of digits stored is sufficient to prevent truncation errors during conversion.

The presentation of date and time is modifiable by users so they can display this information in the appropriate time zone (e.g., India rather than Zulu) and modify it for other zones as needed. Numeric data is properly aligned according to the particular numerical separators and indicators used in the target language. In addition, if an application allows users to manipulate text, the different forms of tabulation available are modified as needed (e.g., allow the decimal tab to work with commas rather than periods) to accommodate the data formats used in the target language.

Number systems and formats. While Arabic numerals (e.g., 0, 1, 2, etc.) are widely accepted, some languages have their own number systems. In some cases (e.g., Chinese), the symbols are substitutes for Arabic numerals, while in others (e.g., Ethiopia), there are special characters for numbers such as 10 and 100.

When presenting numbers, a comma, period, space, and apostrophe can be used as separators for units of thousands. In some cases, an explicit separator is not required for numbers less than 10,000. Numbers can be grouped by thousands or ten thousands. The period, comma, and center dot can be used as separators for decimal numbers. Positive and negative numbers can be indicated by + and - symbols appearing either before or after the number, and negative numbers can be enclosed in parentheses.

Measurement systems and arithmetic operations. US users are familiar with the Imperial system of measurement that uses inches and fractions of an inch (e.g., halves, quarters, eighths) while users outside the US rely on the metric system which measures in meters, liters, and grams. In addition, the US relies on the Fahrenheit scale for temperature while the rest of the world uses Celsius. Similarly, cultures vary in the manner in which certain arithmetic operations are performed. For example, some countries have rules for rounding numbers that differ from those used in the US. In addition, accounting rules (e.g., to calculate compound interest) vary from locale to locale.

Currency. The comma, period, and colon can be used as separators for currency. Currency indicators include a number of symbols (e.g., \$, British pound, and the Japanese yen), alphabetic characters (e.g., FF, SFRs, kr), and combinations (e.g., CZ\$), and can be placed at the beginning, middle, or end of the currency expression. There can be one or no space between the currency symbol and the amount, and currency symbols can be up to four characters in length. Most currencies (except Japan) include two digits to indicate fractional money amounts.

Date and time. The hyphen, comma, period, space, and slash can be used as separators for the day, month, and year, or separators can be left out altogether. In numeric date formats, the month and day fields can be reversed, and in some cases, the year field can come first. Month and day names can be capitalized or in lower case and can be abbreviated using the first two or three letters or some other combination of letters.

The manner in which a date is expressed can be affected by the calendar system being used. While dates are usually based on the Gregorian calendar, some cultures use lunar calendars or the Jewish or Arabic calendar or can express the date based on the year of accession of the Emperor, as in Japan. These calendars can include day names for more than seven days, and month names for more than twelve months. Moslem countries such as Saudi Arabia and Egypt use a calendar with 12 months but only 354 or 355 days. The first day of the week is Sunday in the US but Monday in European countries, a difference that affects the manner in which calendars are displayed.

The colon, period, and space can be used as separators for hours, minutes, and seconds. The letter h can separate hours and minutes. Both 12-hour and 24-hour notation can be used. For 12-hour notation, a.m. or p.m. can appear after the time.

Although the world is divided into 24 standard time zones, countries have the freedom to set their own times. For example, in South America, Surinam's time is 30 minutes different from that of the next zone, and Guyana's is 45 minutes different. The same time zone can have multiple names, and different time zones can share the same abbreviation. Finally, countries differ in their rules concerning daylight savings time or may not use it at all, and the hemispheres differ in when it starts and ends because the seasons are reversed.

Addresses and telephone numbers. Addresses vary from two to six lines long and can include any character used in the character set for a language. The house number precedes the street name in the US and United Kingdom but follows the street name in most other European countries. Postal codes appear in various positions and can include alphabetic characters (e.g., an abbreviation for the country), separators (usually spaces), and numbers (up to seven characters and numbers in length). In many countries, each part of an address is written on a separate line; however, in South Korea, the entire address is placed on a single line, with the specific format used varying for central cities and local areas.

Telephone numbers can contain blanks, commas, hyphens, periods, and square brackets as separators. Telephone numbers can be displayed in local, national, and international formats. Local formats vary widely. National formats can have an area code in parentheses, while international formats can drop the parentheses but add a plus sign at the beginning of the number to indicate the country code.

## **17.4.4 Graphics**

Icons and symbols. The icons and symbols used by an application may be unfamiliar to users outside the US. For example, a mail application that changes a mailbox graphic to indicate receipt of new mail may be unrecognizable in another culture where mailboxes have a different appearance or may not be used. Certain images, colors, and numbers of objects in a group may evoke a negative reaction in another culture so they obscure or contradict the message they are intended to convey. As a result, the icons used in an application may need to be modified in order to match the image or symbol to the culture in which the application will be used.

Whenever possible, an application uses international symbols in its icons. If a new symbol is created, it represents a basic, concrete concept because concrete icons require less explanation than abstract ones. In addition, each new symbol needs to be compared with existing symbols to ensure there are no conflicts. The use of stars and crosses as part of the symbol is avoided. Text is not included in an icon graphic because it will need to be translated and may not fit into the icon when presented in the target language.

Drawings. An application incorporates translated text and adjusts data formats as needed when presenting graphic information (e.g., line graphs, bar charts, flow charts). The size of the graphics objects may need to be enlarged to accommodate the increased length of translated text. Alternatively, application graphics can be modified to place text adjacent to, rather than within, the object so changes in text length do not affect size of individual objects or the overall illustration.

Graphic design conventions vary from culture to culture. For example, Japanese artists tend to draw tables of data differently than Western artists do. As a result, an application may require modification to accommodate these conventions.

Tactical symbology. When an application presents tactical data (e.g., in a map window), users are able to access a variety of map features in order to customize the display to match their preferred mode for viewing and interpreting this information. For example, where US users are likely to display road features for navigation in urban areas, Korean operators may prefer to see neighborhood names as key map landmarks.

Visual cues for alerting. The specific visual signals used by an application, especially for alerting, are reviewed to ensure that they convey the desired meaning in the target culture and that their representations within the software are not objectionable to users. Alert and warning messages can be supplemented with icons so an application communicates critical information in both text and graphic form.

### **17.4.5 Keyboard Interaction**

Mnemonics and shortcut keys. When menu options are translated, any mnemonics or shortcut keys included with the options need to be modified to reflect the translated text. In general, the guidelines for mnemonics in languages with single-byte character sets also apply to languages with multi-byte character sets, except for how mnemonics are displayed. An application translated from the former to the latter can retain the mnemonics used in the single-byte version, with the mnemonic displayed in parentheses following the text of the menu option. If all of the characters in a menu option have been assigned as mnemonics or if the choice consists of multi-byte characters, an application can use another letter or keyboard character. The same mnemonic is assigned to an option whenever it appears in a menu.

The layout of the user's keyboard needs to be considered when selecting the key combinations for the mnemonics and shortcut keys to assign to translated menu options. First, keys are selected to minimize the disruption or relearning required to execute a mnemonic or shortcut key, especially a frequently used one. Second, there are no conflicts between the key combinations for entering accented characters (e.g., if a US keyboard is being used) and those being used for mnemonics and shortcut keys. Finally, some non-US keyboards contain only one ALT, located either on the left or right side of the keyboard. The ease with which users can execute the key combination for a mnemonic is considered if one of these keyboards is being used.

Speed search and text search. If an application is translated into a language that contains accents on the first letter of words, users are able to perform a speed search (see section 3.2.2.3) by typing an unaccented upper-case or lower-case letter, and the search finds instances of both unaccented and accented first letters.

If an application uses wild card characters to perform text searches (see section 12.1.5), it needs to determine if these characters are assigned special meaning in the target language. If necessary, an alternate set of wild card characters is selected to eliminate any possible confusability when an application is converted to the target language.

### 17.4.6 Text Manipulation

Sorting and collation. An application makes use of linguistic sort sequences to order the contents of alphanumeric lists or to add new information to an already sorted list. In US software, characters are usually compared according to their binary value in the code set, with characters ordered on the basis of these values. However, variations are frequently required to reflect linguistic conventions since the binary sequence of characters may not match the linguistic sequence for the language. For example, variations may be needed to handle characters with functional equivalence (e.g., Mac and Mc usually appear together) and to address situations where a character should be ignored (e.g., re-locate and relocate should be placed together). In addition, where US software typically provides a single sorting algorithm to accommodate such variations, other languages usually support multiple sort orders. As a result, an application needs to provide users with the ability to choose a sort order that meets their needs.

Sorting rules for European languages must be able to handle extended character sets and language-specific conventions, independent of the binary values assigned to characters. These languages may contain letters after “z” or sort letters out of the standard alphabetic sequence used in the US. For example, some of these languages contain double characters that sort as one combined character, or a single character that is treated as a double character. In Spanish, double characters such as “ch” and “ll” sort as a single character, and in German, ß is a single character that is treated as “ss” when found in a word. Madell, Parsons, and Abegg in Developing and Localizing International Software provide the following examples of differences in sorting order based on ASCII and German rules, and ASCII and Spanish rules:

| <u>Sorted by<br/>ASCII rules</u> | <u>Sorted by<br/>German rules</u> | <u>Sorted by<br/>ASCII rules</u> | <u>Sorted by<br/>Spanish rules</u> |
|----------------------------------|-----------------------------------|----------------------------------|------------------------------------|
| Airplane                         | Airplane                          | chaleco                          | cuna                               |
| Zebra                            | ähnlich                           | cuna                             | chaleco                            |
| bird                             | bird                              | día                              | día                                |
| car                              | car                               | llave                            | loro                               |
| ähnlich                          | Zebra                             | loro                             | llave                              |
|                                  |                                   | maíz                             | maíz                               |

In the case of complex (i.e., multi-byte) languages, expressions can be written in a mixture of character sets. For example, the Japanese word for “water” may be written as a single Kanji character, as two Hiragana characters, as two Katakana characters, or as the four-letter Romaji expression “mizu.” As a result, sorting algorithms in these languages must be able to accept multiple character patterns as representing the same expression. These algorithms can combine a sorting order among the character sets with a sorting order for expressions within each set. In addition, an application may need to provide a sort order based on a symbol feature that is not captured within the character code. For example, Chinese expressions may need to be sorted by the numeric value of the character as represented in the coded character set as well as by the number of strokes required to represent the character, the radical (i.e., root) of the character, or the number of strokes added to the radical. Finally, an application may need to implement a sort order based on the way symbols are pronounced. In this case, each symbol may have to be stored in both graphic and phonetic form, with the resulting sort order listing symbols that are phonetically similar but visually different near each other.

Editing functions. Editing functions (e.g., search and replace, cut and paste, and spell checking) can accommodate the unique features of the target language, including instances where the appearance of a word changes when it is hyphenated, where it appears in lower rather than upper case, or where it contains a combination character such as æ. In contextual languages such as Thai, the characters that make up a compound symbol may be entered in several different orders, with the appearance of the symbol dependent on the order in which the characters are entered. In other languages (e.g., Greek), the appearance of a character can vary depending on its position in a word. If an application performs string searches in these languages, it is able to recognize any of several possible character sequences and judge them to be the same or different as appropriate.

### **17.4.7 Adjustments for Bi-directional Languages**

If an application is localized to a bi-directional language such as Hebrew or Arabic, window orientation and information orientation within each window are adjusted as appropriate for right-to-left presentation. Window appearance is the mirror image of that in English-based windows, except that the location of the window buttons in the title bar does not change. In addition, window placement is oriented right-to-left; i.e., a primary window is positioned to the right and its child window(s) to the left.

With respect to information orientation in bi-directional languages, an application complies with the style specifications in this document, except that “right” and “left” are interchanged. However, physical right and left remain the same. As in unidirectional languages, LEFT and RIGHT move the cursor in the arrow direction; the right and left buttons on the pointing device behave as defined here, with left and right movement of the pointing device moving the pointer in these directions.

With respect to information content, an application provides translations for titles, headings, prompts, and other window controls, except for English acronyms not normally

translated and the names of keys on the keyboard. If an application chooses to mix right-to-left and left-to-right elements within the same window, it follows the relevant specifications defining information display for unidirectional and bi-directional languages.

### **17.4.8 Adjustments for Vertical Languages**

Asian languages such as Chinese and Japanese contain a combination of horizontally and vertically written characters, with the latter written from top to bottom and each new line starting to the left of the previous one. If an application is converted to a vertical language, it provides translations for titles, headings, prompts, and other window controls, except for English acronyms not normally translated and the names of keys on the keyboard. An application displays this information in the orientation expected by users, complying with relevant specifications defining information display for horizontal and vertical languages. In addition, an application supports data entry for vertically written text using one of the text-input methods described in section 17.3 and presenting a vertically-oriented text entry area.

### **17.4.9 Printing**

Peripheral devices such as printers are capable of handling the character set for the target language; i.e., the full character set can be loaded on the printer, and the printer can produce all of the extended characters required by the language.

While the standard paper size in the US is 8.5 x 11 inches, most countries use ISO A4 size which is slightly longer and narrower than the US standard. As a result, printer capabilities (e.g., different paper trays) may need to be adjusted in order to handle the standard paper and envelope sizes used by the target audience, and an application may need to be modified to handle the varying page layouts dictated by the different paper sizes. For example, hardcoded rules regarding paper margins are removed, and users are allowed to specify how they want text to appear and to do so using measurement units with which they are familiar.

Adjustments made in window format to accommodate text expansion also need to consider text presentation when the content of the window is printed. In particular, the amount of vertical space between lines of text is sufficient to print all extended characters, including those with accents, in both upper and lower case.

### **17.4.10 Internationalizing Web Applications**

While it is desirable to present the content of a Web application in a user's preferred language, the cost to translate the entire application into each of the languages used by its audience can be prohibitive. An alternative is to provide translations for some pages and leave others in the original language. This hybrid approach requires that an application decide on a default language, with users able to switch between available languages on the home page and lower-level pages as needed. If an application cannot decide on a

default language for its home page, it can provide a staging page where users select their preferred language before navigating to the home page. Whenever users are presented with a choice of languages, the list should contain the name of the language as a word (rather than as a graphic such as a national flag), using the language's own name for itself.

Because it is unlikely that all of the content of a Web application will be translated into every language it supports, the application needs to support a multilingual search capability in order to cover the entire information space of the application. This capability should be designed so that users enter the desired search terms in their preferred language and the application translates the terms into the requested languages before performing the search. This approach is preferred to having users identify the appropriate translated synonyms in each of the languages, which may result in inaccurate or incomplete searches.